# Supplementary figures

**Contents**

**Supplementary figure legends**

**Figure S1 Evaluating the mutation rate of the primers.** The upper part shows the frequency distribution of the four bases. The middle represents the major allele. The lower part shows the depth of random NGS sequencing. More than 98% virus RNA contained the identical sequences with the amplification primers of AGGGGTACTGCTGTTATGTCT and CCTTCAGATTTGTTCGCGC.

**Figure S2 Details of the top 30 abundant mutants of each sample.** Each line represented a kind of mutant, green lines were the missense mutations and the red colors were synonymous mutations. The top rectangle represents the cumulative variations from quasispecies mutants. Functional domain of the spike gene were marked with different colors: NTD, N-terminal domain; RBD, receptor-binding domain; S1/S2, protease cleavage site; FP, fusion peptide; HR1, heptad repeat 1; CH, central helix; CD, connector domain; HR2, heptad repeat 2; TM, transmembrane domain; CT, cytoplasmic tail. The rightmost bar plots showed the proportion of mutants (abundance). Results of haplotype clustering were listed at the right bottom, where mutation spectrums centered by the master one were identified.

**Figure S3 The occurrence statistics of the consensus genomes with 3118G>T in Genbank database (by March 2021).** The results showed that 3118G>T was never emerged before October 2020 that was 10 month after the sampling date of the current study.

**Figure S4 The position of 3118G>T on the spike protein (3D ribbon structure).** The corresponding amino acid is highlighted by a red line and some potential functional domains are marked with arrow. The nucleotide of 3118T was located in the turn site between CH and CD sub-domains.

**Figure S5 Simulation for the amino acid changes of the 3118 G>T.** Part **a** showed the nucleotide of SG master mutant G (same with the early ancestor strain), and part **b** showed the nucleotide of FG master mutant T. The location of variation site was highlighted by yellow color, the changed angle of R-group was simulated by Chimera with the Posterior probability=0.30.

**Figure S6 Principal component analysis for all samples.** (**a**) Two clusters were identified based on all SNVs of the top 30 abundant mutants. (**b**) The SNV 3118T>G as the first dimension could fully distinguish the two single-source infection groups and made more than 98% contributions to the grouping.
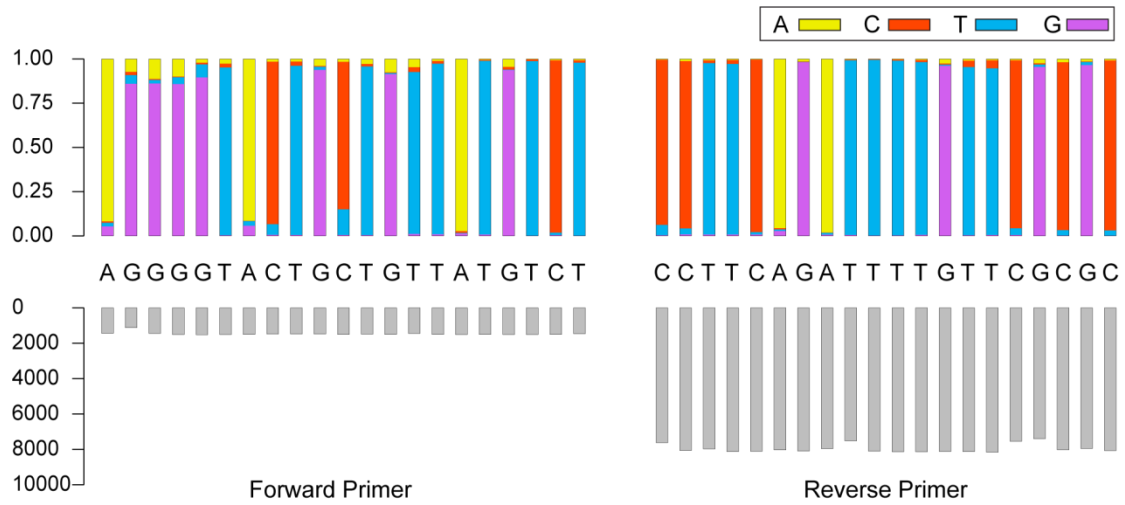
**Figure S7 Occurrence of three hub spike haplotypes of Genbank in the quasispecies mutants.** Three haplotypes reflected for three master mutants within patients, and all these master mutants had already existed in the minor mutants of early infected patients.

**Figure S8 Detecting a currently predominant SNV 1841A>G (D614G) in the minor mutants of the quasispecies.** The nucleotide 1841G (amino acid G614) had already existed in 19 minor mutants from 4 early infected patients.

**Figure S9 Box plots of the nonsense and missense mutation rate for three sampling sites.** Both percentage of nonsense and missense mutation showed obviously higher in samples of stools. The mean values were marked on the figure.

**Figure S10 Box plots of the dN and dS value.** The mean dS values of each domain are almost the same (**a**), while the mean dN values of three domains of FP, HR1, and CD were significantly larger than the whole gene level (t-test) (**b**).

FG1-0126-SP: Average Major Allele Frequency 0.9454±0.0382

FG1-0129-NS: Average Major Allele Frequency 1.00

FG2-0127-NS: Average Major Allele Frequency 0.99±0.01

FG4-0127-NS: Average Major Allele Frequency 0.9994±0.0039

SG2-0131-SP: Average Major Allele Frequency 0.9910±0.02758

SG3-0205-ST: Average Major Allele Frequency 0.95±0.03

**Figure S1 Evaluating the mutation rate of the primers.** The upper part shows the frequency distribution of the four bases. The middle represents the major allele. The lower part shows the depth of random NGS sequencing. More than 98% virus RNA contained the identical sequences with the amplification primers of AGGGGTACTGCTGTTATGTCT and CCTTCAGATTTGTTCGCGC.

Top 30 Abundant Mutants of SG1-0131-SP

Top 30 Abundant Mutants of SG2-0131-SP

Top 30 Abundant Mutants of SG4-0211-ST

Top 30 Abundant Mutants of SG3-0203-NS

Accumulated Mutations in MN908947.3

Functional Domains

NTD    RBD    S1/S2    FP    HR1    CH    CD

Relative Abundance

Equal to Reference ( An early strain from Wuhan: MN908947.3 )    61.97%

G1337T:G446V    1.65%

C3407T:T1136I    0.7%

A3763T:K1255U    0.64%

T1129C:F377L    0.56%

A3292G:N1098D    0.52%

G2049T:R683R    0.51%

G2872T:A958S    0.51%

G3371T:G1124V    0.51%

G3792T:V1264V    0.5%

A3215G:E1072G    0.48%

A2905G:N969D    0.46%

G1223T:R408I    0.45%

G1087A:A363T    0.45%

T3647C:I1216T    0.45%

G1966T:V656F    0.42%

G1943T:G648V    0.41%

C3424T:Q1142U    0.4%

A912G:K304K    0.39%

T2609G:I870S    0.38%

C3026T:T1009I    0.37%

T48C:V16V    0.36%

A2351G:Q784R    0.33%

T1978C:Y660H    0.32%

G489A:A163A    0.31%

C2293T:R765C    0.31%

G2006C:G669A    0.31%

G456T:W152C    0.29%

G1943T:G648V    G2054T:R685L    0.27%

C35T:S12F    0.27%

Master mutant
61.97%

○ Relative abundance (Area size)
● Mutants with 3118T genotype
● Mutants with 3118G genotype

■ Synonymous mutation    ■ Missense mutation

10

Top 30 Abundant Mutants of SG3-0205-ST

Top 30 Abundant Mutants of SG3-0206-SP

Accumulated Mutations in MN908947.3

Functional Domains

NTD     RBD     S1/S2     FP     HR1     CH     CD

Relative Abundance

Equal to Reference ( An early strain from Wuhan: MN908947.3 )

74.07%

| Mutation | Abundance |
|---|---|
| G3711C:M1237I | 0.14% |
| C145T:H49Y | 0.14% |
| C2102T:A701V | 0.11% |
| T2742C:N914N | 0.1% |
| C3622T:Q1208U | 0.09% |
| C2293T:R765C | 0.09% |
| T2255C:L752P | 0.07% |
| T2813C:L938P | 0.07% |
| T2874C:A958A | 0.07% |
| C623T:T208M | 0.07% |
| G2031T:Q677H | 0.07% |
| T1978C:Y660H | 0.07% |
| G1239T:G413G | 0.06% |
| T2400A:F800L | 0.06% |
| C1640T:T547I | 0.06% |
| T823C:F275L | 0.06% |
| T3019C:Y1007H | 0.06% |
| T2910G:F970L | 0.06% |
| C3759T:C1253C | 0.06% |
| T1773C:S591S | 0.05% |
| T3657G:G1219G | 0.05% |
| T1128C:T376T | 0.05% |
| T1129C:F377L | 0.05% |
| T2298C:A766A | 0.05% |
| T3675C:I1225I | 0.05% |
| G1855T:E619U | 0.05% |
| G3296A:G1099D | 0.05% |
| C2636T:A879V | 0.05% |
| G224A:G75D | 0.05% |

Master mutant
74.07%

○ Relative abundance (Area size)
● Mutants with 3118T genotype
● Mutants with 3118G genotype

■ Synonymous mutation     ■ Missense mutation

12

Top 30 Abundant Mutants of FG1-0126-NS

Top 30 Abundant Mutants of FG1-0126-SP

Top 30 Abundant Mutants of FG1-0127-NS

Top 30 Abundant Mutants of FG1-0129-NS

Top 30 Abundant Mutants of FG2-0127-NS

# Top 30 Abundant Mutants of FG3-0127-NS



Accumulated Mutations in MN908947.3

Functional Domains

NTD  RBD  S1/S2  FP  HR1  CH  CD

Relative Abundance

| Mutation labels | Relative Abundance |
|---|---|
| G3118T:V1040F | 69.99% |
| G3118T:V1040F / G3371T:G1124V | 1.61% |
| G24C:L8F / G3118T:V1040F | 1.45% |
| G3118T:V1040F / T3177G:G1059G | 0.9% |
| G3118T:V1040F / T3463A:Y1155N | 0.8% |
| C2648T:T883I / G3118T:V1040F | 0.77% |
| C2521T:L841F / G3118T:V1040F | 0.73% |
| G1781T:G594V / G3118T:V1040F | 0.73% |
| T1284C:D428D / G3118T:V1040F | 0.72% |
| G3118T:V1040F / T3129C:C1043C | 0.71% |
| G1337T:G446V / G3118T:V1040F | 0.69% |
| A2734G:T912A / G3118T:V1040F | 0.67% |
| G1943T:G648V / G3118T:V1040F | 0.67% |
| A606T:K202N / G3118T:V1040F | 0.56% |
| C547T:Q183U / G3118T:V1040F | 0.52% |
| T3096G:C1032W / G3118T:V1040F | 0.49% |
| T2157C:T719T / G3118T:V1040F | 0.49% |
| C2121A:Y707U / T2188C:S730P / G3118T:V1040F | 0.48% |
| T47C:V16A / G3118T:V1040F | 0.47% |
| C1425T:A475A / G3118T:V1040F | 0.41% |
| T1945G:C649G / G3118T:V1040F | 0.37% |
| T1522C:Y508H / G3118T:V1040F | 0.37% |
| A2432G:K811R / G3118T:V1040F | 0.36% |
| G3118T:V1040F / G3624T:Q1208H | 0.33% |
| T47C:V16A / G663A:S221S / G3118T:V1040F | 0.3% |
| C1425T:A475A / A2328T:K776N / G3118T:V1040F | 0.29% |
| C905T:T302M / G3118T:V1040F | 0.24% |
| G3118T:V1040F / C3299T:T1100I | 0.22% |
| A21G:L7L / G3118T:V1040F | 0.16% |
| T495G:N165K / A546T:K182N / G3118T:V1040F | 0.14% |

Synonymous mutation    Missense mutation

Master mutant 69.99%

○ Relative abundance (Area size)
● Mutants with 3118T genotype
● Mutants with 3118G genotype

18

Top 30 Abundant Mutants of FG4-0127-NS

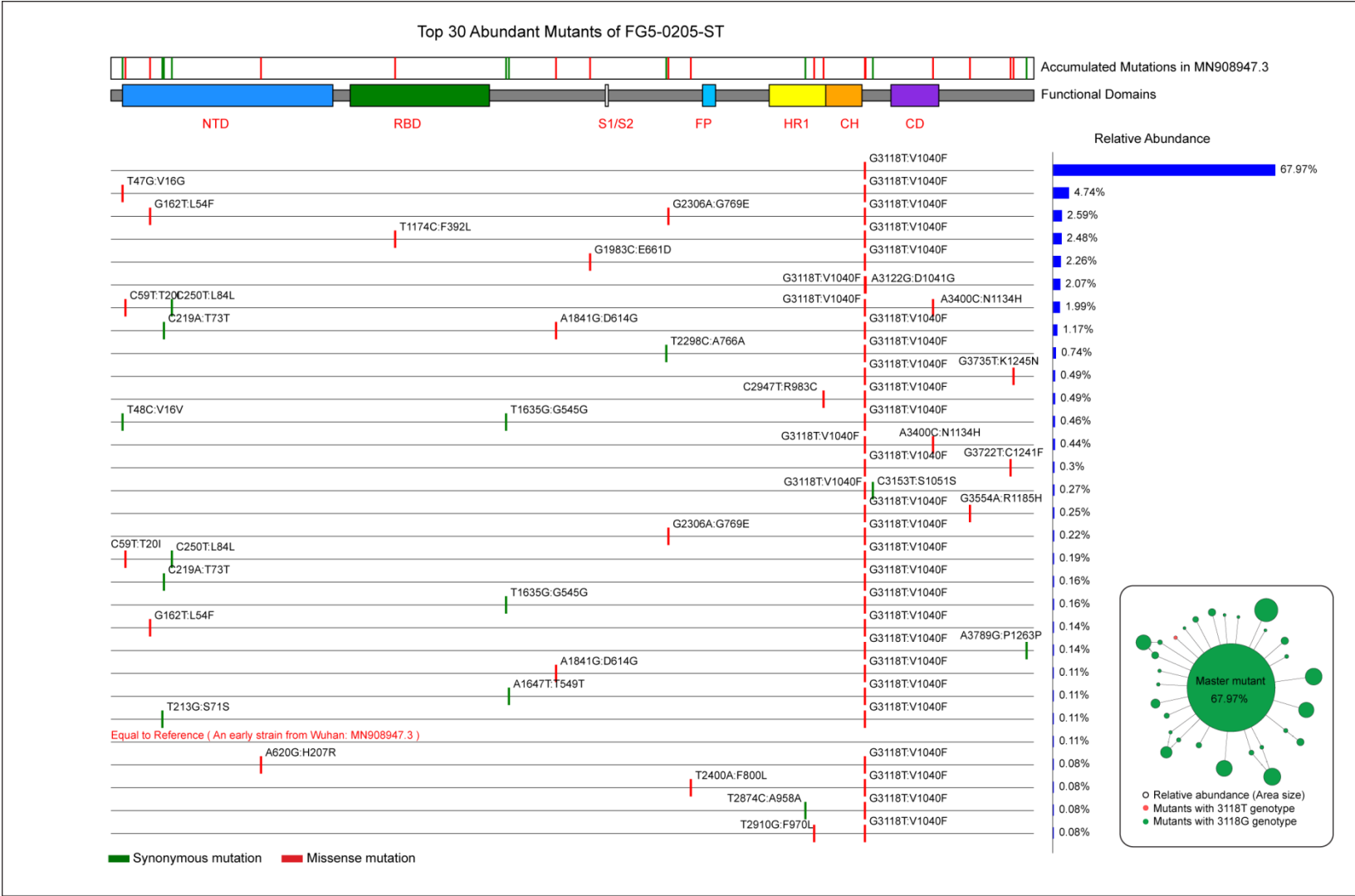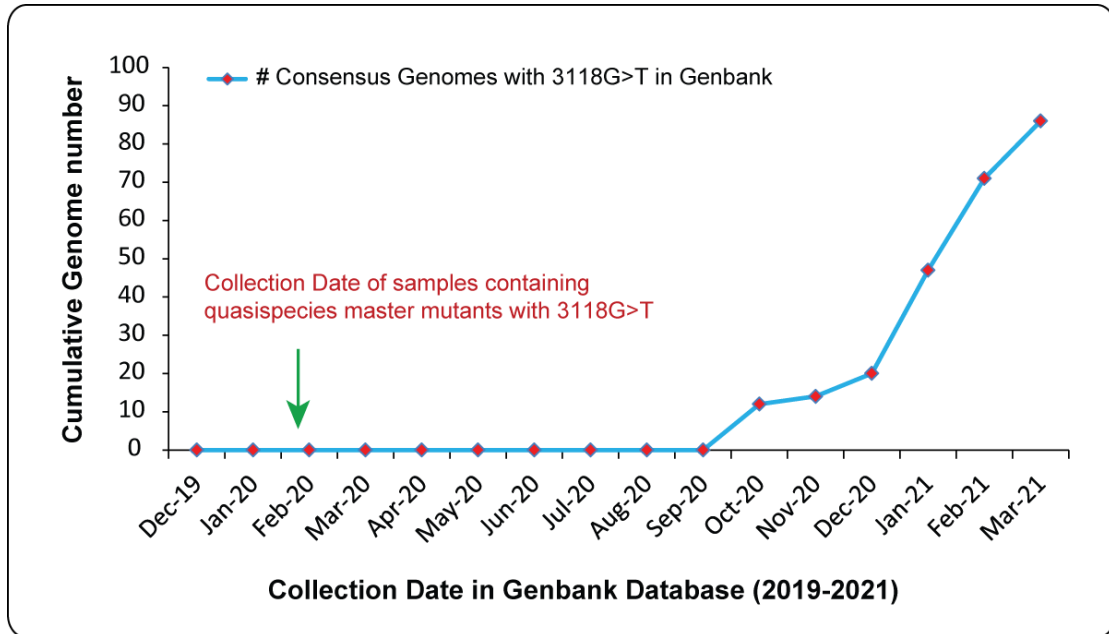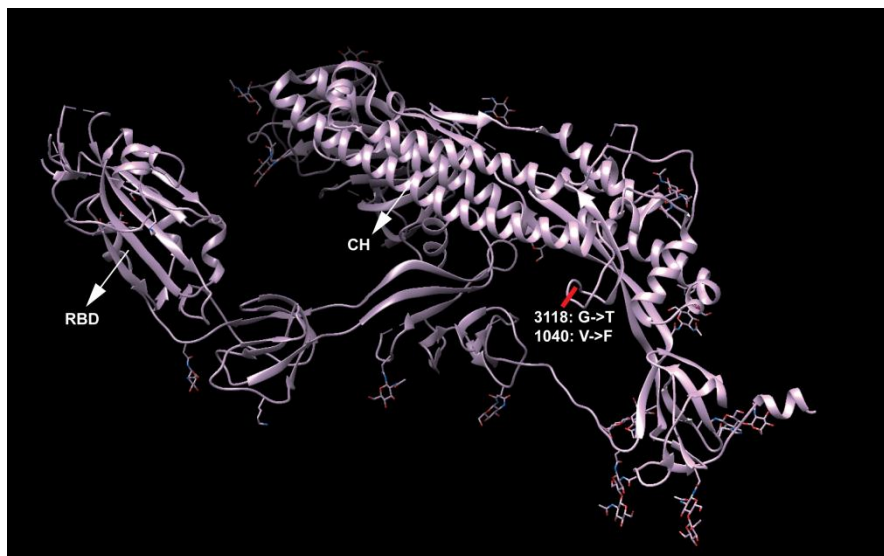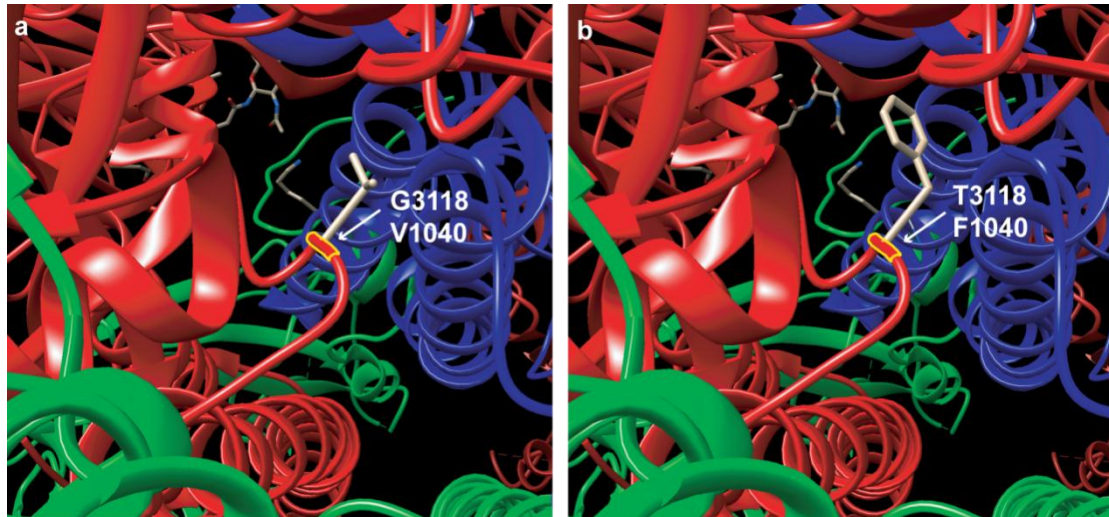Top 30 Abundant Mutants of FG5-0205-ST

**Figure S2 Details of the top 30 abundant mutants of each sample.** Each line represented a kind of mutant, green lines were the missense mutations and the red colors were synonymous mutations. The top rectangle represents the cumulative variations from quasispecies mutants. Functional domain of the spike gene were marked with different colors: NTD, N-terminal domain; RBD, receptor-binding domain; S1/S2, protease cleavage site; FP, fusion peptide; HR1, heptad repeat 1; CH, central helix; CD, connector domain; HR2, heptad repeat 2; TM, transmembrane domain; CT, cytoplasmic tail. The rightmost bar plots showed the proportion of mutants (abundance). Results of haplotype clustering were listed at the right bottom, where mutation spectrums centered by the master one were identified.
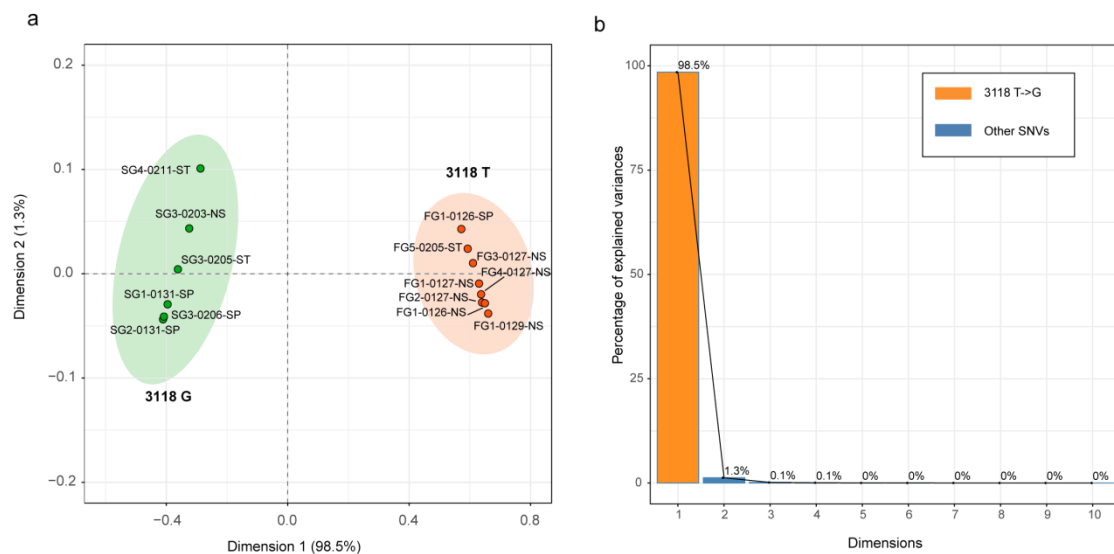
**Figure S3 The occurrence statistics of the consensus genomes with 3118G>T in Genbank database (by March 2021).** The results showed that 3118G>T was never emerged before October 2020 that was 10 month after the sampling date of the current study.
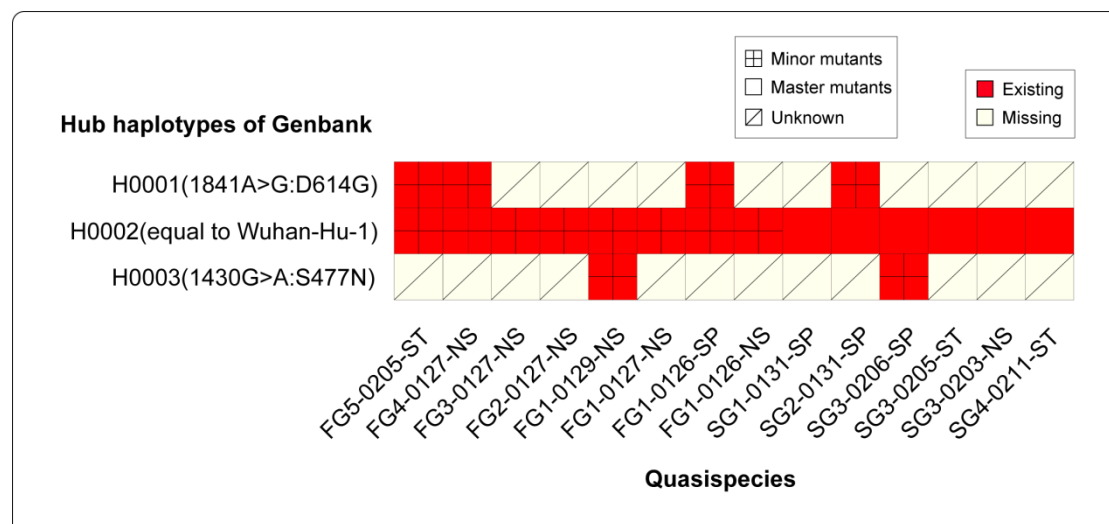


**Figure S4 The position of 3118G>T on the spike protein (3D ribbon structure).** The corresponding amino acid is highlighted by a red line and some potential functional domains are marked with arrow. The nucleotide of 3118T was located in the turn site between CH and CD sub-domains.
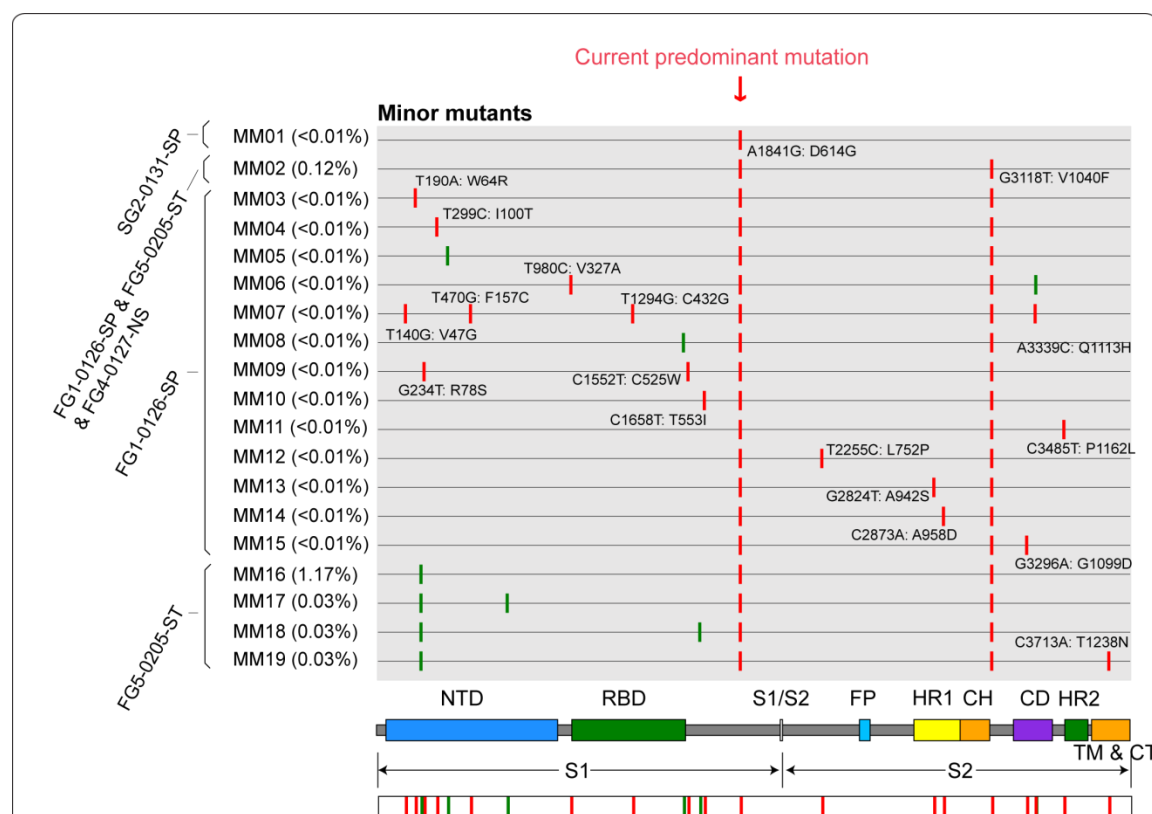
**Figure S5 Simulation for the amino acid changes of the 3118 G>T.** Part **a** showed the nucleotide of SG master mutant G (same with the early ancestor strain), and part **b** showed the nucleotide of FG master mutant T. The location of variation site was highlighted by yellow color, the changed angle of R-group was simulated by Chimera with the Posterior probability=0.30.



**Figure S6 Principal component analysis for all samples.** (**a**) Two clusters were identified based on all SNVs of the top 30 abundant mutants. (**b**) The SNV 3118T>G as the first dimension could fully distinguish the two single-source infection groups and made more than 98% contributions to the grouping.
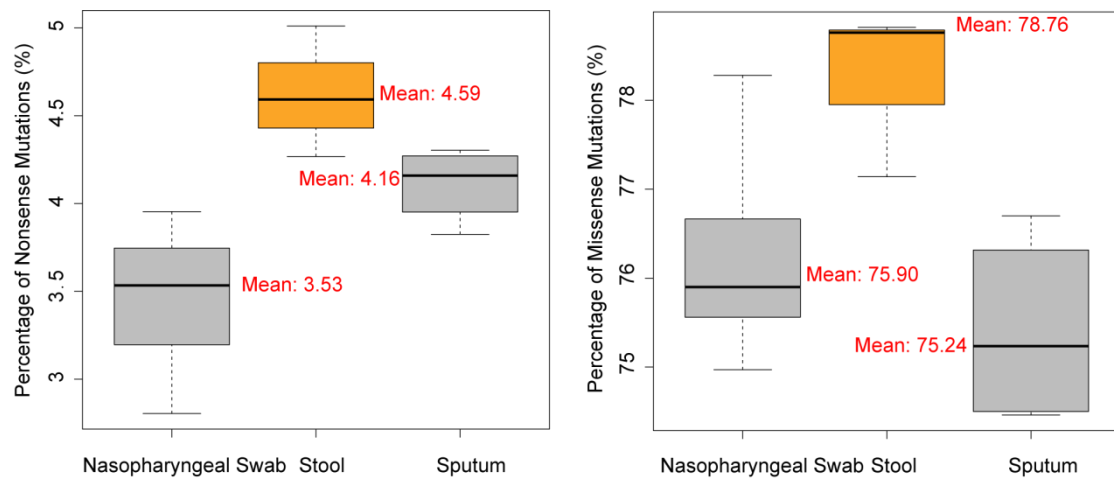
**Figure S7 Occurrence of three hub spike haplotypes of Genbank in the quasispecies**

**mutants.** Three haplotypes reflected for three master mutants within patients, and all these

master mutants had already existed in the minor mutants of early infected patients.
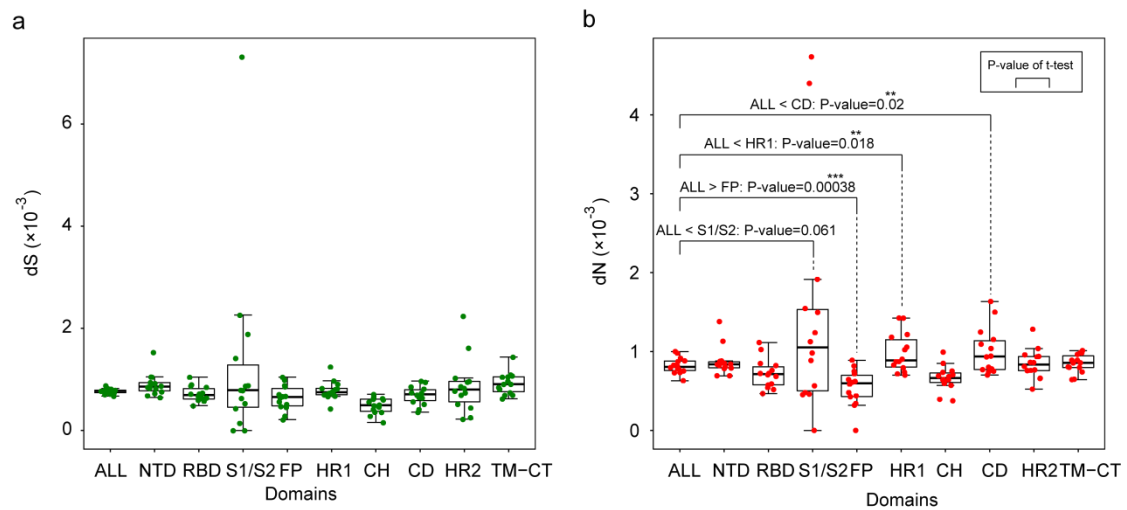


**Figure S8 Detecting a currently predominant SNV 1841A>G (D614G) in the minor**

**mutants of the quasispecies.** The nucleotide 1841G (amino acid G614) had already existed

in 19 minor mutants from 4 early infected patients.

**Figure S9 Box plots of the nonsense and missense mutation rate for three sampling sites.**
Both percentage of nonsense and missense mutation showed obviously higher in samples of stools. The mean values were marked on the figure.



**Figure S10 Box plots of the dN and dS value.** The mean dS values of each domain are almost the same (**a**), while the mean dN values of three domains of FP, HR1, and CD were significantly larger than the whole gene level (t-test) (**b**).